

A load balancing device according to an embodiment of the invention uses a predictor that comprises a plurality of Least Connections Control Blocks (LCCBs) that keeps track of the real servers with active connections. To speed up the search for the real server with the least number of active connections, an LCCB is kept for each metric. A metric is defined as the number of connections on a server divided by its weight (or capacity) of the server. This metric is kept as a quotient/remainder pair. The predictor sends out the real server address with the lowest metric whenever a new connection is required by the load balancing device.